# Automatic Phonetic Transcription by Phonological Derivation

Marcos Garcia[1] and Isaac J. González[2]

[1] Center for Research in Information Technologies (CITIUS)
University of Santiago de Compostela
[2] Cilenis Language Technology
**Draft Version**
marcos.garcia.gonzalez@usc.es,isaacjgonzalez@cilenis.com
http://gramatica.usc.es/pln
http://www.cilenis.com

**Abstract.** Automatic phonetic transcription tools usually perform phonetic transcriptions directly from orthographic representations. Although these approaches often achieve good results, theoretical studies suggest that including morphophonological knowledge allows those systems to improve their performance. Following this idea, we developed a tool which first obtains an underlying representation of each word, using small lexica and dedicated lemmatizers. For each representation, a phonological derivation generates the phonetic transcription by applying linguistically motivated rules. Since most of these rules are added as optional parameters, the system permits to generate dialect-specific transcriptions. This system is not only a grapheme-to-phone tool, but it also obtains phonological representations and evaluates several linguistic processes occurring during the derivation. Preliminary experiments emulating a phonological system of Galician (using as input words spelled in European Portuguese) show that the underlying representation of most words can be obtained using small lexica and also that the derivation produces high-quality phonetic transcriptions.

**Keywords:** grapheme-to-phoneme, phonetics, phonology, galician, portuguese

## 1 Introduction

Automatic Phonetic Transcription (APT) is a crucial task for many applications of different areas. Besides text-to-speech systems, which need high quality transcriptions, APT tools are also used in theoretical and applied linguistics. These tools are useful in many areas (e.g., phonetics, phonology, dialectology, language learning, etc.) in order to obtain preliminary transcriptions of large corpora.

Rule-based APT systems often generate a phonetic transcription directly from the orthographic form. These approaches achieve good performance in languages whose spelling systems permit to *easily* infer their phonetic representation, such as Italian, Spanish or Portuguese. However, algorithms which do not

take into account morphophonological information may produce some errors due to their lack of this linguistic knowledge [1]. In this respect, the transcription of Galician and Portuguese varieties presents some problems. For example, the orthographic vowels <e> and <o> could represent up to six phones (e.g., [ɛ, e, ɨ, o, ɔ, u]), some of them not being predictable from their context.

The transcription of these not predictable examples could be performed through rules and lists of exception words. Nevertheless, if there is no morphophonological knowledge, these lists may be very large, since they should include inflected forms: nominal forms (from 'festa') such as 'f[ɛ]sta', 'f[ɛ]stas', 'f[ɛ]stinha' (or 'f[e]stinha') and verbal forms (from 'levar') such as 'l[ɛ]vo', 'l[ɛ]vas', 'l[ɛ]va', etc. Thus, the creation of these lists would be time-consuming.

In this paper, a strategy to overcome these limitations is proposed. The behaviour of the phonemes in the root of the words is consistent and predictable among their inflected forms, so we use shorter lists of exceptions, including only those lemmas which have characters which do not follow the default conversion rules. This way, by applying a rule-based lemmatizer on the inflected words —as well as a set of specific rules for each linguistic variety— the system obtains the target lemma and verifies whether it is in an exception list. This process allows the system to generate an underlying (phonological) representation, which is an abstract form of the word before the application of the phonological rules.

Then, a phonological derivation produces phonetic transcriptions by applying rules on the underlying representations, as shown in Figure 1.

| <festinha> / <festiña> | Orthographic Representation |
|---|---|
| <festa> (Exception: root vowel: /ɛ/) | *Lemma* |
| **/fɛstiɲa/** | **Underlying Representation** |
| /fɛs.ti.ɲa/ | *Syllable Split* |
| /fɛsˈti.ɲa/ | *Stress assignment* |
| /fɛsˈtĩ.ɲa/ | *Nasalization* |
| /fɛsˈtĩ.ɲɐ/ | *Unstressed vocalism* |
| **[fɛsˈtĩ.ɲɐ]** | **Phonetic Representation** |

**Fig. 1.** Example of the conversion of the word 'festinha/festiña'.

It is worth noting that different rules —belonging to several dialects— generate diverse transcriptions from the same underlying form. Thus, the selection of the rules involves different outputs (for dialects sharing the same underlying form). In our system, dialectal rules were added as optional parameters, so it is capable of automatically transcribing lexica for different dialects.

The system presented in this paper (released under GPL license,[3] and whose first version was presented in [12]) emulates the phonological system of Galician varieties by following the mentioned strategy. The input can be written in two different spellings: ILG/RAG [14] and European Portuguese (EP).

---

[3] `www.gnu.org/licenses/gpl.txt`

The main advantages of using the proposed method are that the user can (i) obtain phonological representations, (ii) analyze phonological processes such as syllabification, (iii) evaluate the phonological derivation as well as interact with it, and (iv) create phonetic representations for different linguistic varieties.

Experiments performed on a Portuguese corpus show that the system obtains underlying representations accurately. Furthermore, the results also indicate that the phonological derivation generates high quality phonetic transcriptions.

This paper is organized as follows: Section 2 shows some related work concerning automatic phonetic transcription of Galician and Portuguese. In Section 3, we briefly introduce the theoretical background of our method, whose architecture is presented in Section 4. Then, Section 5 contains the performed experiments and, finally, conclusions and further work are addressed in Section 6.

## 2 Related Work

Much of the work on automatic phonetic transcription has been done focused on text to speech synthesis. However, there is also other approaches related to our work that have to be taken into account.

Rule-based models were implemented in order to perform grapheme-to-phone conversion of Galician (and EP) varieties, achieving results of more than 98% precision [3, 5]. Another work focused on the automatic transcription of Galician presents the main characteristics and difficulties of this task [13]. Besides the segmental features mentioned above (the transcription of <e>, <o> or <x>), this paper also introduces some other aspects, such as contraction forms or intonation issues. Moreover, the development of a corpus and a lexicon for Galician text-to-speech systems is presented in [8].

The disambiguation of heterophonic homographs is another important process of the automatic phonetic transcription task. For Galician, the main difficulties are presented in [19].

There are much work focused on the APT of Portuguese varieties, from rule-based to stochastic models [4, 21, 25, 27]. For our goals, it is interesting to refer [26], which improves the precision of syllable splitting by applying phonological theories, namely the onset-rhyme theory. Another work which includes phonological processing is [24], which generates phonetic variants for speech recognition.

In [1], is presented a project whose aim is the population —taken into account morphological information— of a large lexicon of Portuguese with different phonetic transcriptions of several linguistic varieties.

Finally, another tool which uses linguistic knowledge is FreP [28]. FreP automatically extracts data about frequency and linguistic contexts of phonological entities, allowing a phonologist to easily obtain this kind of information.

## 3 Theoretical Background

The system presented in this paper was designed following some phonological theories, that we briefly describe in this section.

The main idea of our method is the use of several representation levels as proposed in classical generative phonology [9]. This theory describes the phonology of languages as derivational systems with different abstraction levels. In this view, the underlying form of a word —also known as *deep* or *phonological* representation— is transformed into a surface (or phonetic) form by the application of transformational rules in the phonological derivation process.

Lexical phonology [18], proposed as a refinement of classical generative phonology, introduced two main components in the description of the phonology of a language: lexical and post-lexical. In the first one, information of different levels (morphology, syntax, etc.) plays a role in the derivation, while the post-lexical component has no access to that knowledge. The application of post-lexical rules on the output of the lexical component generates the final surface form.

From this point of view, it can be postulated that the members of a linguistic community share the underlying forms of their language. So, different rules (or rule order) will generate several surface forms, corresponding to different variants of the same phonological word. For instance, we can assume that the underlying form of 'feira' for EP speakers is /feira/. Thus, different realizations such as 'f[ej]ra', 'f[ɐj]ra' or 'f[e]ra' are generated due to differences in the derivation.

Classical generative phonology did not take into account syllables as units of analysis, so the explanation of many phonological processes within syllables were rarely simple. Different non-linear models of phonology emerged in order to introduce the syllable —and its internal structure— into phonological theory.

Onset-rhyme (syllable) theory defines syllables as phonological units which organize segmental melodies in terms of sonority [2]. Moreover, in this theory syllables have an internal structure, exemplified in Figure 2. The phonology of a particular language will define (i) if a syllabic constituent (Onset, Coda) is obligatory or optional, (ii) what segments could occupy each position, as well as (iii) the maximum number of segments which can be anchored to each constituent. This structure is useful to explain, for instance, some phonological processes which affect phonemes depending on the syllabic position in which they occur.

These (and other) theories were applied for describing the phonology of Portuguese varieties [17]. Taking into account the differences, much of this work is useful for analyzing Galician dialects. Specifically for Galician, some phonological and morphological aspects are described in works such as [7, 22, 10].
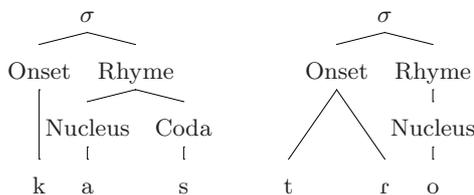


**Fig. 2.** Internal structure of the syllables for the word 'castro' (Onset-Rhyme Theory).

The phonological theories described in this section were taken into account in the development of our system. They were implemented with minor variations, in order to build an APT tool capable of generating different phonetic representations from a single underlying form. Furthermore, note that the system is also useful to obtain deep phonological representations of words as well as to automatically verify the behaviour of derivational rules on large corpora.

In this paper, we use the terms 'phoneme' and 'phone' to designate phonological and phonetic entities, respectively.

## 4   Architecture

Our system is composed of several modules applied in a pipeline. The different components are described in this section.

### 4.1   Phonological Conversion

The first step of our method is the conversion of an orthographic word into an underlying representation. Two different modules perform this process, depending on the spelling of the input (ILG/RAG or EP orthographies).

Both modules are compound of a set of transformational rules as well as lists of exceptions. The rules simply substitute orthographic characters by representations of phonological segments taking into account their context (<v>→ /b/, <ch>→ /t͡ʃ/, <rr>→ /r/, <c(aou)>→ /k(aou)/, etc.).

Those characters whose representation is not predictable by their context are included in the exception lists. We use the following, extracted from [23] (and verbal forms automatically extracted from large lexica):

– List of lemmas with <x> as /ks/ ('se/ks/o').
– Lists of lemmas with <e> as /ɛ/ and with <o> as /ɔ/ ('f/ɛ/sta', 'r/ɔ/da').[4]

By default, <x>, <e> and <o> are transformed into /ʃ/, /e/ and /o/. The words included in these lists are spelled in ILG/RAG, so we use a transliteration method in order to perform a look-up with words written in EP. For this spelling, we also use the following lists, semi-automatically extracted:

– List of lemmas with <qu/gu> followed by <e/i> as /kw/, /gw/ (fre/kw/ente).
– List of lemmas ending in <ão> as /an/ (p/an/).
– List of exceptions: *non-galician* forms ('fiz' → /fiʃen/, 'sim' → /si/).

We have to note that the transliteration strategy may produce some errors, since some words are not easily transformed with this method. However, its performance was successfully tested in several cases [15, 16].

---

[4] To be more precise, lists of words with /ɛ, ɔ/ are not pure lemmas, but singular forms containing these vowels in the root. This way, the first list contains forms such as the adjective 'nova', but not the plural 'novas'.

Apart from the referred lists, two sets of Galician heterophonic homographs (one for each spelling system) were compiled, each one of about 600 pairs. The representation of these words depends on their PoS-tag (verb: '/ɔ/lho'; noun: '/o/lho'), semantic meaning (*weapon*: 'b/ɛ/sta'; *animal*: 'b/e/sta') or verbal mood (present: 'c/ɔ/me'; imperative: 'c/o/me'). The system does not perform disambiguation, so it generates both outputs of the heterophonic homographs.

**Dedicated Lemmatizer:** The use of large lists of exceptions (including all the inflected forms which behave as their lemmas) is avoided as follows:

An open-source rule-based lemmatizer which detects if the input word is an inflected form is applied. Then, the system verifies whether the input (or the obtained lemma) exists in the exception list. If it exists, the underlying representation is generated by applying the exception rules of each list (<e>→ /ɛ/, <x>→ /ks/, etc.). The lemmatizer is compound of two sub-modules: one for nominal forms (inspired in [6]), and a dedicated stemmer for verbs. Figure 3 shows some examples of this process.

Nominal metaphony in Galician dialects is usually different from Portuguese metaphony. Galician plurals maintain the vowel of the singular form ('[o]lho / [o]lhos', and not '[o]lho / [ɔ]lhos'), but the vowel changes in both singular and plural feminine forms ('n[ɔ]va / n[ɔ]vas': the masculine singular of these adjectives could be 'n[o]vo' or 'n[ɔ]vo' depending on the dialect). This way, and taking into account that the lists of these words also include singular feminine forms, the nominal lemmatizer only applies rules to infer the singular form.

Note that this implementation slightly differs from lexical phonology. In our system, some morphophonological processes (such as nominal and verbal inflection) are replaced by the use of lists and by lemmatization. Furthermore, the underlying forms could be different to others proposed in theoretical literature, which may include morpheme information or other phonological abstractions.

The output of the phonological conversion module is an underlying representation of the input word, obtained by including morphophonological processing. This output is the input of the next module, which performs syllabification.

### 4.2   Syllabification and Stress Assignment

The function of this module is to split the input word into syllables and to build the internal structure of each syllable, following the onset-rhyme theory. Then, the stressed syllable is detected and marked.

| Input | | Lemma | | Exception? | | Underlying |
|---|---|---|---|---|---|---|
| <levam> | $\xrightarrow{lemmatizer}$ | levar | $\xrightarrow{look-up}$ | yes | $\xrightarrow{Conversion\ rules}$ | /lɛban/ |
| <festas> | $\xrightarrow{lemmatizer}$ | festa | $\xrightarrow{look-up}$ | yes | $\xrightarrow{Conversion\ rules}$ | /fɛstas/ |
| <cestas> | $\xrightarrow{lemmatizer}$ | cesta | $\xrightarrow{look-up}$ | no | $\xrightarrow{Conversion\ rules}$ | /sestas/ |

**Fig. 3.** Examples of the underlying form generation.

A simplified variant of the algorithm proposed in [17], adapted to the Galician syllable structure, was implemented. The sonority hierarchy was removed, so we established which segments (and groups of segments) could occupy each syllabic position. Thus, in order to build the syllable structure, we follow these steps:

1. Falling diphthongs are marked as nucleus (/p<u>ei</u>ʃe/).
2. Vowels are marked as nucleus (/peiʃ<u>e</u>/).
3. Complex onsets are marked (/<u>pr</u>atos/).
4. Simple onsets followed by nucleus are marked (/pra<u>t</u>os/).
5. Complex (and *non-patrimonial*) codas are marked (/a<u>bs</u>traiɾ/).
6. Simple codas are marked (/abstrai<u>ɾ</u>/).
7. Exception (uncommon codas, foreign words, etc.) are identified (/ga<u>ngs</u>ter/).

Then, syllable boundaries are marked (i) before each onset, (ii) before an initial nucleus, (iii) between two nucleus and (iv) between a coda and a nucleus.

Once syllable split is performed, the stressed syllable of each word has to be marked, in order to apply the phonological rules accurately. The system maintains the orthographical input during the conversion, so this process applies the accentuation rules inversely. It first verifies if the word has an accent, marking as stressed the syllable with it. If there is no accent, the inverse accentuation rules are applied until find the stressed syllable. There is also a small set of unstressed words, such as prepositions and pronouns ('com', 'me', 'te', etc.).

### 4.3 Derivational Rules

Phonological rules are applied after the syllabification process and the stress assignment, allowing them to incorporate syllabic knowledge. However, it is worth noting that some derivational rules may imply re-syllabification processes.

Derivational rules are applied sequentially, so the order of application is crucial for obtaining the desired output. Current version of our system contains both a set of *universal* rules (occurring in most of Galician dialects), and optional rules (corresponding to specific phonological processes which do not occur in every Galician variety), all of them present in the theoretical literature. The selection of some optional rules allows the user to obtain transcriptions of the standard variety as well as of non-standard Galician dialects [23].

- Universal rules:
  - Semi-vocalization: this rule transforms some adjacent vocalic phonemes into a diphthong (/so.si'al/ → [so'sjal]).
  - Unstressed vocalism: unstressed vowels are elevated and centralized in some contexts, namely in final word position (/'ko.mo/ → ['ko.mʊ]).
  - Voiced plosives: voiced plosives /b, d, g/ are pronounced as approximants ([β̞, ð̞, ɣ̞]) in some contexts (/a'bɾiɾ/ → [a'β̞ɾiɾ]).
  - Nasalization: this rule changes the place of articulation of implosive nasals. It could perform an assimilation (the nasal segment acquires the place of articulation of the following consonant: /'kan.pʊ/ → ['kam.pʊ])

or a velarization (implosive nasals are always velar: [ˈkaŋ.pʊ]). Further-more, an additional rule nasalizes the vowels occurring in some contexts, such as between nasal segments (/ˈmaŋ.ta/ → [ˈmãŋ.ta]).

- Optional rules:
  - Thetacismo: this is a special rule to deal with a main characteristic of Galician dialects. Some of them have these two sibilants phones: [s, θ] ('caçar': [kaˈθaɾ]; 'casar': [kaˈsaɾ]), while others only have [s] ('caçar', 'casar': [kaˈsaɾ]). This exceptional rule is not derivational, and it also af-fects the underlying representation of the word, including both phonemes /s, θ/ or only /s/, depending on the application.
  - *Gheada*: another dialect-specific process. This rule changes every /g/ by [ħ], except when occurring after a nasal segment, allowing some other realizations ('água': [ˈa.ħwa]; 'bingo': [ˈbiŋ.gʊ], [ˈbiŋ.ħʊ] or [ˈbiŋ.kʊ]).
  - Complex codas: this rule simplifies complex codas as well as removes some *non-patrimonial* codas (/aβsˈtra<u>k</u>.tʊ/ → [asˈtra.tʊ]).
  - Rhotacism: this rule affects the phoneme /s/ in some coda positions (depending on the features of the following consonant), changing its re-alization by [ɾ] (/ˈmes.mʊ/ → [ˈmeɾ.mʊ]).
  - Voicing: this rule will represent /s, θ/ in coda (when occurring before voiced onsets) by their voiced allophones (/ˈdes.dɪ/ → [ˈdez.dɪ]).

This module implements —with some differences— the behaviour of a phono-logical system as proposed by the theories described in Section 3. Rules are similar to those of classic generative phonology, but taking advantage of the knowledge provided by onset-rhyme theory as well as of regular expressions. This way, the addition, modification and deletion of the rules is a simple task that allows the user to test the behaviour of new phonological rules.

Note that the system performs phonological and phonetic conversion of in-dividual words, so it does not include post-lexical rules applying across word boundaries. Moreover, current version of the system neither has a specific mod-ule for proper nouns, numeric expressions and acronyms.

## 5   Experiments

We carried out several evaluations in order to know the performance of the system transcribing into Galician a corpus written in European Portuguese.[5]

For this purpose, we randomly selected from CETEMPúblico a 10,000 token corpus.[6] Abbreviations, numbers and acronyms without vowels were removed.

First, the corpus was tokenized and PoS-tagged with FreeLing [20, 11], al-lowing us to disambiguate several heterophonic homographs. Then, the system automatically transcribed each word. The chosen Galician variety —selected through the optional rules— was those proposed in [23] (with minor changes, due to the use of a different spelling). This output was manually revised by an

---

[5] A previous version of this system was evaluated on a ILG/RAG corpus in [12].
[6] http://www.linguateca.pt/cetempublico/

**Table 1.** Types of error of the orthography to phonology conversion (*Full* evaluation).

| Type of Error | Errors |
|---|---|
| <x> | 8 |
| <e/o> | 168 |
| <qu/gu> | 1 |
| <ão/ões> | 2 |
| *foreign words* | 63 |

expert, creating a gold standard corpus. For those words for which [23] offers more than one choice, the gold standard only includes the first transcription.[7]

The obtained corpus has about 8,500 tokens (without punctuation), 18,000 syllables and 40,000 phonemes. The internal structure of the syllables were not revised, so the evaluation of the syllabification process only takes into account the boundaries between these elements.[8]

Evaluations of two processes were performed: "orthography to phonology" and "phonology to phonetics". In order to carry out these evaluations, the errors produced during the conversion were manually revised, classifying them according to the two mentioned processes (note that PoS-tagging errors were computed as errors of the system). Furthermore, two different results were calculated: *Full*, taking into account the transcription of the whole corpus and *noForeign*, which does not evaluate the transcription of proper nouns and foreign words.

In the evaluation of the "orthography to phonology" module, we counted the errors produced during the grapheme-to-phoneme conversion, thus ignoring rising diphthongs (Table 1). Most of the errors concern the transcriptions of <e> and <o>, mainly due to errors in the transliteration of the input. Furthermore, eight of them were produced by PoS-tagging errors.

In order to evaluate the "phonology to phonetics" conversion, we computed the errors produced by (i) the syllabification, (ii) the stress assignment as well as (iii) the derivational modules.

First, the precision of the syllabification was calculated, also taking into account the derivational rules which involve re-syllabification. The system produced 97 errors (75 with a *noForeign* evaluation) on 14,155 syllable boundaries (the splits between words were not evaluated), achieving a precision of 99.31%.

In our system, stress assignment is highly influenced by the previous process (syllable splitting). Thus, 62 of the 102 errors produced in this stage were motivated by previous splitting errors. Nevertheless, these 102 errors point out that this module has a precision of 99.46%.

Table 2 contains the results of the two processes. The first two lines show that the orthography to phonology conversion achieved a precision of more than 99%.

---

[7] The software and the gold standard corpus will be available at the following website: `http://gramatica.usc.es/~marcos/software/`

[8] Syllable boundaries were annotated following the formal register described in [22]. In addition, rising diphthongs were considered as phonetic rather than phonological: <social>: /sɔθial/ → /sɔ.θi.al/ → [...] → [sɔˈθjal].

**Table 2.** Results of the two processes as well as the full phonetic conversion.

| Process | Evaluation | Phonemes | Errors | Precision |
|---|---|---|---|---|
| Orthography → Phonology | *Full* | 39,394 | 242 | 99.39% |
| | *noForeign* | 36,052 | 153 | 99.58% |
| Phonology → Phonetics | *Full* | 39,394 | 199 | 99.48% |
| | *noForeign* | 36,052 | 130 | 99.64% |
| Orthography → Phonetics | *Full* | 39,394 | 441 | 98.88% |
| | *noForeign* | 36,052 | 283 | 99.22% |

Central lines indicate the precision of the phonology to phonetics conversion, which also scored more than 99%. Finally, the bottom lines of Table 2 show the results of the whole conversion (orthography to phonetics).

A deep evaluation of the errors shows that, as said above, many of them were produced by orthographic issues. The use of external resources (such as the exception lists and the lemmatizers) in a different spelling system necessarily involves some errors during the conversion. Other errors were generated by not listed exceptions (namely cases of hiatus/diphthong, and few errors produced by the lemmatization strategy), by foreign/uncommon words and by heterophonic homographs whose disambiguation needs semantic information. However, taking the above into account, the obtained results show that the performance of our method is comparable to other state-of-the-art systems for similar languages.

## 6    Conclusions and Further Work

In this paper, an open-source system which performs automatic phonetic transcription by phonological derivation was presented.

The performed experiments, although preliminary, show that the use of morphophonological processing permits to obtain accurate phonological representations (even from two different spelling systems). Moreover, the application of phonological theories as well as the use of linguistically motivated rules overcome some limitations of previous grapheme-to-phoneme approaches, allowing the system to automatically generate high-quality phonetic transcriptions of different dialects.

Further work will be focused on correcting and increasing the exception lists, on the inclusion of dedicated rules for guessing the vowel height of unknown words as well as on the improvement of the transliteration strategy. Moreover, the system also needs to include a preprocessing module capable of analyzing abbreviations, numbers and acronyms.

Finally, we plan to evaluate this method on the automatic transcription of European and Brazilian Portuguese dialects, by taking into account the morphophonological differences compared to Galician varieties.

# References

1. Ashby, S., Ferreira, J. P.: The Role of Morphology in Generating High-Quality Pronunciation Lexica for Regional Variants of Portuguese. In: Proceedings of the 9th International Conference on Computational Processing of the Portuguese Language (PROPOR'10). LNAI, vol. 6001, pp. 162–165. Springer-Verlag (2010)
2. Blevins, J.: The Syllable in Phonological Theory. In: Goldsmith, J. A. (ed.), The Handbook of Phonological Theory, pp 206–244. Blackwell, Cambridge (1995)
3. Braga, D., Coelho, L.: Letter-to-sound conversion for Galician TTS systems. In: Actas de las IV Jornadas en Tecnologia del Habla, pp. 171–176. Zaragoza (2006)
4. Braga, D., Coelho, L., Resende Jr., F.: A Rule-Based Grapheme-to-Phone Converter for TTS Systems in European Portuguese. In: Proceedings of the VI International Telecommunications Symposium (ITS'06), pp. 328–333. Fortaleza (2006)
5. Braga, D., Freixeiro, X. R.: Algoritmos de Conversão Grafema-Fone em Galego para Sistemas de Conversão Texto-Fala. In: Estudos galegos de Tradución & Paratradución no século XXI. Xerais, Vigo (2007)
6. Branco, A., Silva, J.: Very High Accuracy Rule-Based Nominal Lemmatization with a Minimal Lexicon. In: Actas do XXI Encontro Anual da Associação Portuguesa de Linguística (2007)
7. Castro, O.: Aproximación a la fonología y morfología gallegas. PhD Thesis, Georgetown University (1989)
8. Campillo, F., Braga, D., Mourín, A. B., García-Mateo, C., Silva, P., Sales Dias, M., Méndez F.: Building High Quality Databases for Minority Languages such as Galician. In: Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC'10). ELRA, La Valleta (2010)
9. Chomsky, N., Halle, M.: The Sound Pattern of English. Harper and Row, New York (1968)
10. Dubert García, F.: Máis sobre o rotacismo de /s/ en galego. In: Álvarez, R., Vilavedra, D. (eds.), Cinguidos por unha arela común. Homenaxe ó profesor Xesús Alonso Montero, pp. 367–387. Universidade de Santiago de Compostela (1999)
11. Garcia, M., Gamallo, P.: Análise Morfossintáctica para Português Europeu e Galego: Problemas, Soluções e Avaliação. Linguamática. Revista para o Processamento Automático das Línguas Ibéricas 2(2), 59–67 (2010)
12. Garcia, M., González, Isaac J.: Conversión Fonética Automática con Información Fonológica para el Gallego. Procesamiento del Lenguaje Natural 47, 283–291 (2011)
13. González González, M., Banga, E. R., Campillo, F., Méndez, F., Rodríguez Liñares, L., Iglesias, G.: Specific features of the Galician language and implications for speech technology development. Speech Communication 50, 874–887 (2008)
14. ILG/RAG: Normas Ortográficas e Morfolóxicas do Idioma Galego. Real Academia Galega and Instituto da Lingua Galega, Vigo (2005)
15. Malvar, P., Pichel, J. R., Senra, Ó., Gamallo, P., García, A.: Vencendo a escassez de recursos computacionais. Carvalho: Tradutor Automático Estatístico Inglês-Galego a partir do corpus paralelo Europarl Inglês-Português. Linguamática. Revista para o Processamento Automático das Línguas Ibéricas 2(2), 31–38 (2010)
16. Malvar, P., Pichel, J.R.: Generación semiautomática de recursos de Opinion Mining para el gallego a partir del portugués y el español. In: ICL: Workshop on Iberian Cross-Language NLP tasks. 27th Conference of the Spanish Society for Natural Language Processing. Huelva (2011)
17. Mira Mateus, M. H., Andrade, E. d'.: The Phonology of Portuguese. Oxford University Press, Oxford (2000)

18.  Mohanan, K. P.: The Theory of Lexical Phonology. Dordrecht, Reidel (1986)
19.  Mourín, A., Braga, D., Coelho, L., García-Mateo, C., Campillo, F., Dias, M.: Homograph Disambiguation in Galician TTS Systems. In: IX Congreso Internacional da Asociación Internacional de Estudos Galegos. A Coruña - Santiago de Compostela - Vigo (2009)
20.  Padró, Ll.: Analizadores Multilingües en FreeLing. Linguamática. Revista para o Processamento Automático das Línguas Ibéricas 3(2), 13–20 (2011)
21.  Paulo, S., Oliveira, L. C., Mendes, C., Figueira, L., Cassaca, R., Viana, C., Moniz, H.: DIXI — A Generic Text-to-Speech System for European Portuguese. In: Proceedings of the 8th International Conference on Computational Processing of the Portuguese Language (PROPOR'08). LNAI, vol. 5190, pp. 91–100. Springer-Verlag (2008)
22.  Regueira, X. L.: A sílaba en galego: lingua, estándar e ideoloxía. In: Lorenzo, R. (ed.), Homenaxe a Fernando R. Tato Plaza, pp. 235–254. Universidade de Santiago de Compostela, Santiago de Compostela (2002)
23.  Regueira, X. L.: Dicionario de Pronuncia da Lingua Galega. Real Academia Galega and Instituto da Lingua Galega, A Coruña (2010)
24.  Seara, I. C., Pacheco, F. S., Seara Júnior, R., Kafka, S. G., Klein, S., Seara, R.: Geração Automática de Variantes de Léxicos do Português Brasileiro para Sistemas de Reconhecimento de Fala. In: Actas do XX Simpósio Brasileiro de Telecomunicações. Rio de Janeiro (2003)
25.  Siravenha, A. C., Neto, N., Macedo, Vl, Klautau, A.: Uso de Regras Fonológicas com Determinação de Vogal Tônica para Conversão Grafema-Fone em Português Brasileiro. In: Proceedings of the 7th International Information and Telecommunication Technologies Symposium (I2TS'08). Foz do Iguaçu (2008)
26.  Oliveira, C., Castro Moutinho, L., Teixeira, A. J. S.: On European Portuguese automatic syllabification. In: Proceedings of Interspeech 2005, pp. 2933–2936 (2005)
27.  Veiga, A., Candeias, S., Perdigão, F.: Generating a Pronunciation Dictionary for European Portuguese Using a Joint-Sequence Model with Embedded Stress Assignment. In: Proceedings of the 8th Brazilian Symposium in Information and Human Language Technology (STIL'11), pp. 144–153 (2011)
28.  Vigário, M., Martins, F., Frota, S.: A ferramenta FreP e a frequência de tipos silábicos e classes de segmentos no Português. In: Selecção de Comunicações apresentadas no XX Encontro Nacional da Associação Portuguesa de Linguística, pp. 675–687 (2006)